



## Commentary

## The aggregation challenge

Macartan Humphreys<sup>a,\*</sup>, Alexandra Scacco<sup>b,\*</sup><sup>a</sup> Columbia University and Director, IPI Unit, WZB Berlin, Germany<sup>b</sup> Senior Research Fellow, IPI Unit, WZB Berlin, Germany

## ARTICLE INFO

## Article history:

## ABSTRACT

Banerjee, Duflo, and Kremer have had an enormous impact on scholarship on the political economy of development. But as RCTs have become more central in this field, political scientists have struggled to draw implications from proliferating micro-level studies for longstanding macro-level problems. We describe these challenges and point to recent innovations to help address them.

© 2019 Elsevier Ltd. All rights reserved.

Banerjee, Duflo, and Kremer have had an enormous impact on work on the political economy of development. Although there were parallel trends in political science (e.g., [Green & Gerber, 2002](#)), much of the early expansion of RCTs into topics in political economy was not just inspired by but led by the Nobel laureates and their close associates. Key studies on leadership, aid politics, and accountability include [Chattopadhyay and Duflo \(2004\)](#), [Gugerty and Kremer \(2008\)](#), and [Banerjee, Kumar, Pande, and Su \(2011\)](#). [Olken \(2007\)](#) and [Miguel, Satyanath, and Sergenti \(2004\)](#) pioneered the use of experiments and natural experiments to address questions about corruption and political violence. Today, hundreds of researchers are implementing experiments to study the political economy of development. There are well over a thousand registered experimental designs on the Evidence in Governance and Politics (EGAP) registry alone. Experiments have taken root.

At the same time, there are broad concerns that the experimental turn has diverted researchers from core questions in our field. Indeed, there has been an awareness of this risk from the outset. In [2006](#), Robert Bates wrote that “Banerjee’s approach might teach us more about impact but at the expense of larger matters,” warning political scientists against transforming the field “from a search for the underlying forces of development into a form of policy analysis.” Despite this concern, attempts by political scientists to aggregate the lessons learned from rigorous micro-level experimental work to shed light on larger puzzles have been casual at best.

## 1. Micro-macro disconnects

The aggregation problem is distinct from the problem of the external validity of experimental results, though that’s a part of

it. Rather, the problem is like trying to figure out how pieces of a jigsaw puzzle fit together when many of the pieces are missing.

Three examples of aggregation problems:

- You are interested in whether freedom of the press fosters better government. What can you learn from an experiment that shows that voters are more likely to vote against politicians when they learn that the politicians are underperforming?
- You are interested in whether larger endowments of natural resources weaken state-society linkages. What can you learn from an experiment that shows that voters who are told that revenues are derived from natural resources – rather than taxes – exhibit less concern about government expenditures?
- You are interested in whether inter-ethnic violence is caused by residential segregation. What can you learn from an experiment that shows that prejudice decreases among individuals exposed to higher levels of contact with out-group members?

In all three cases, a macro-level question motivates the research, but the researcher has micro-level experimental evidence at hand. The micro-level experiments seem to provide relevant evidence, but it is unclear what inferences to draw from this micro evidence for the macro questions.

Let’s take the third example and use it to flesh out three distinct parts of the aggregation problem. We are interested in whether ethnic conflict is exacerbated by segregated settlement patterns. A macro-level hypothesis (drawn from several rich, largely non-experimental literatures in political science) might be that certain residential patterns – such as ethnically mixed but highly segregated cities – can heighten feelings of insecurity which, when exploited by opportunistic political elites, play an important role in explaining the incidence of communal violence.

[Fig. 1](#) illustrates with a simple graphical model, linking residential segregation and conflict at the macro level. Different parts of

\* Corresponding authors.

E-mail addresses: [macartan.humphreys@wzb.eu](mailto:macartan.humphreys@wzb.eu) (M. Humphreys), [alex.scacco@wzb.eu](mailto:alex.scacco@wzb.eu) (A. Scacco).

the model find support in different bodies of research. Kasara (2017) uses observational data at the locality level to document the relationship between segregation and levels of intergroup violence in Kenya. Allport (1954) work on the negative relationship between contact and prejudice sparked decades of work on this link. Horowitz (1985) has explored connections between prejudice and ethnic conflict in multiple settings, and a rich literature has linked the strategic behavior of political elites to communal violence (Wilkinson, 2004).

There are obvious reasons to be worried about multiple forms of confounding for the macro-level model shown in Fig. 1. But there is little scope, for practical and ethical reasons, for experimentation with macro-level settlement patterns. Even though experimentalists have repeatedly shown that randomization can be used for a much wider set of problems than skeptics expected, we take it as given that major “treatments,” such as conflict histories and demographic structures, are out of reach. Can (individual) micro-level interventions help?

Fortunately, the theoretical accounts in macro-level studies often specify micro-level logics (for example, Kasara interprets her macro findings through the lens of a micro-level hypothesis about interpersonal mistrust in settings of segregation). Recent experimental studies have examined some of these micro hypotheses directly. In one example, Scacco and Warren (2018) conducted an education-based field experiment in which 850 randomly sampled Christian and Muslim young men in a riot-prone city in Nigeria were randomly assigned to religiously mixed or homogeneous computer training courses. After four months of intergroup contact, they found significant declines in discrimination among subjects assigned to mixed classes.

The question is whether and how inferences from such a contact study can help us understand the macro-level relationship between residential segregation and prejudice, or segregation and violence.

Making direct inferences in this example is difficult for a number of reasons. We highlight three.

### 1.1. Limited learning about selection processes

A great advantage of randomization is that it can overcome selection biases. This achievement creates a problem, however, if the macro processes you want to study involve self-selection. Following the example illustrated in Fig. 1, we can imagine situations where, even if values on the macro-level node (say, “the share of individuals interacting with out-group members”) are randomly assigned, the values of micro-level nodes (whether a given individual interacts with out-group members) might not be.

Imagine, for example, that for a given level of societal segregation, the individuals that encounter out-group members self-select from among those for whom social contact has the weakest effect. Say an experimentalist randomly selects individuals that have rarely been exposed to out-group members and experimentally induces exposure. In this case, the estimate of the average effect would be larger than the average effect for the population, and would thus be a poor estimate of the effect of segregation. The question the researcher needs to answer – what is the average, or overall, effect of the existing level of segregation (relative to some benchmark) for those that self-select into it – is generally not addressed in experimental studies.

*We need to understand selection logics in order to map from micro estimates to macro estimands, but, by design, experimental approaches often prevent us from learning about them.*

### 1.2. Micro averages, macro nonlinearities

Experimental approaches often measure the average effect of a binary treatment, and average effects are treated like linear coefficients

when making inferences regarding aggregate effects. If out-group exposure results in one fewer interpersonal conflict, on average, then treating one person will reduce conflicts by 1 and treating 1000 will reduce conflicts by 1000. Yet, while experiments typically measure average effects *given an overall level of exposure*, the overall level of exposure may also matter for the average effect. And variation in the level of exposure (such as the degree of segregation) may be precisely the variation of interest at the macro level. We need to understand the implied nonlinearities for a host of case-level counterfactual estimands (what would be the level of prejudice in a given country if exposure was at level  $x$  rather than level  $y$ ?). Effects may be different at different exposure levels because of spillover or general equilibrium effects. But there are other reasons this might be the case. For instance, the kind of self-selection into contact we described above produces a nonlinear relationship between the share of individuals exposed to out-groups and the share that is prejudiced.<sup>1</sup> One implication is that studies in different sites can produce different answers, even if the underlying causal processes are the same everywhere.

*The average effect of a micro-level treatment can depend on the overall (macro-level) exposure level, which can produce nonlinearities in effects of exposure levels. Understanding these nonlinearities can be important for macro level attribution questions and may require multisite studies.*

### 1.3. Aggregation requires understanding rival pathways

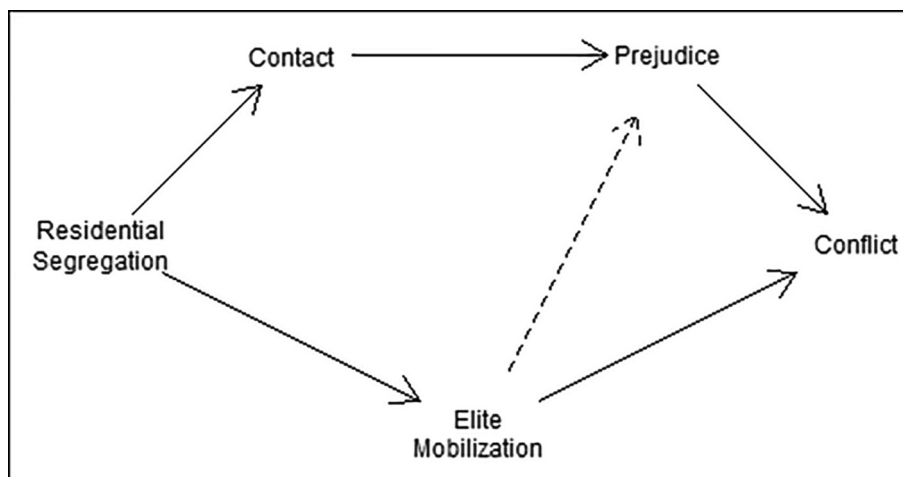
If our goal is to draw inferences about the effects of segregation on intergroup conflict from evidence about the operation of one step in a causal chain – the effects of interpersonal contact on prejudice – then we need to know something about how alternative paths operate. If the relationship between residential segregation and intergroup contact is well understood (for example, residential segregation can reasonably be understood to reduce contact), and if there are no other channels through which segregation affects prejudice, then an argument for aggregation might invoke “the front door criterion” (Pearl, 2009) – where the effect of segregation on prejudice is the effect of segregation on contact multiplied by the effect of contact on prejudice. But if there are rival pathways, then this argument becomes much harder to make, as it is now possible that segregation could increase prejudice, even if there is a negative effect running through individual contact. (See the dotted line in Fig. 1. An example of an alternative path might be via a mechanism of mobilization by opportunistic political elites.)

*We cannot draw implications for the effects of  $X$  on  $Y$  from evidence of effects along links on a causal chain between  $X$  and  $Y$  unless we understand alternative paths from  $X$  to  $Y$ .*

## 2. Pointers toward solutions

These challenges are present even in auspicious settings, where the macro nodes in Fig. 1 are simple aggregates of micro-level measures, where macro variables are as-if-randomly assigned,

<sup>1</sup> A more formal illustration: Imagine a polity with a unit mass of individuals. Say that in the case that share  $k > 0$  individuals encounters out-group members, these, under natural assignment, will be individuals in  $[0, k]$ . Say everyone is prejudiced in the absence of contact, and encountering out-group members eliminates prejudice for individuals in  $[m, 1]$ . Then there is a nonlinear relationship between the share exposed to out-groups and the share that is prejudiced, with no effect for  $k < m$  and a unit marginal effect for  $m < k$ . The effect of moving from no contact to full contact is  $(1 - m)$ . However an experimentalist implementing a contact experiment with individuals randomly selected from among those that have not been exposed to out-groups will find an effect of 1 if  $m < k$  and of  $(1 - m) / (1 - k)$ , otherwise. In both cases the estimate is too large. The actual average treatment effect is  $(1 - m)$ , and the average effect among those treated is  $(1 - m / k)$  for  $m < k$ , and 0 otherwise. Note also that the answer depends upon  $k$ .



**Fig. 1.** Simple model of ethnic residential segregation and conflict linkages. We are interested in the effect of segregation on conflict, both defined at the macro (e.g., national) level. What can we learn from micro level experimental evidence on the effect of contact on prejudice?

and where we assume that experimental estimates correctly estimate effects induced by observational variation.

To make progress addressing the aggregation challenge, we need to grapple with a number of problems that do not figure prominently enough in current research.

First, we need to better understand processes of selection. In order to translate from the individual-level marginal effect to the macro-level marginal effect, we need to understand self-selection rather than simply remove it. One way to do this might be for researchers to randomly partition their study sample into two groups. In one, assignment to treatment is experimentally controlled. In the other, study subjects self-select into treatment. Combining data from these sub-samples would allow for updating on treatment effects, on selection propensities, and on how these relate to each other (Knox, Yamamoto, Baum, & Berinsky, 2019). Another promising approach is provided by Chassang, Miquel, & Snowberg, 2012, in which an experiment simultaneously studies incentives to enter treatment and the effects of treatment, given different propensities to enter treatment.

Second, we need to better understand nonlinearities in relevant relationships at the micro and macro levels. To understand nonlinearities in the macro effect (for example, the relationship between segregation and conflict), we might build on innovations in multi-site experimental studies, as exemplified by the “Metaketa” approach (Dunning et al., 2019). For the contact example, one would want to understand how effects of contact differ in locations with greater or lower levels of baseline exposure to out-group members.

Third, we need to better understand rival pathways. To justify mappings from micro-level evidence to macro-level claims, we need to articulate a *theory* that provides a fuller mapping between macro treatments, micro treatments, and outcomes. Given such a theory, and data on multiple possible channels, we need tools that let us combine inferences in a principled way. The challenges of making causal claims in the presence of multiple channels are formidable (Green, Ha, & Bullock, 2010). Structural approaches provide pointers here, for example, as demonstrated by Pearl and Bareinboim (2014) on transportability and data fusion. At a minimum, this approach requires the researcher to make explicit which relations they take to be general across settings, in order to license combining data on different parts of a causal process.

Finally, we need to actually start doing it. The norm in research in the political economy of development is to generate tight micro-level inferences and then gesture towards macro-level

implications. Doing more to figure out which inferences for larger questions are justifiable will likely require a commitment to articulating how macro conclusions can be justified from micro data, greater re-coordination of research around core substantive agendas, a greater openness to learning from data even when we only enjoy partial identification, and a greater tolerance for deploying models to aid inference – or at least to make explicit the model we already have when we gesture to broader implications of experimental findings.

## References

- Allport, G. W. (1954). *The nature of prejudice*. Addison-Wesley.
- Banerjee, A., Kumar, S., Pande, R., & Su, F. (2011). Do informed voters make better choices? Experimental evidence from urban India. Unpublished manuscript.
- Bates, R. (2006). Banerjee's approach might teach us more about impact but at the expense of larger matters. *Boston Review of Books*, 31(4).
- Chassang, S., Miquel, P. I., & Snowberg, E. (2012). Selective trials: A principal-agent approach to randomized controlled experiments. *American Economic Review*, 102(4), 1279–1309.
- Chattopadhyay, R., & Duflo, E. (2004). Women as policy makers: evidence from a randomized policy experiment in India. *Econometrica*, 72(Sep.), 1409–1443.
- Dunning, T., Grossman, G., Humphreys, M., Hyde, S. D., McIntosh, C., & Nellis, G. (Eds.). (2019). *Information, Accountability, and Cumulative Learning: Lessons from Metaketa I*. Cambridge University Press.
- Green, D., & Gerber, A. (2002). Reclaiming the experimental tradition in political science. In I. Katznelson & H. Milner (Eds.), *Political Science: State of the Discipline* (pp. 805–832). New York: W.W. Norton.
- Green, D. P., Ha, S. E., & Bullock, J. G. (2010). Enough already about “black box” experiments: Studying mediation is more difficult than most scholars suppose. *The Annals of the American Academy of Political and Social Science*, 628(1), 200–208.
- Gugerty, M. K., & Kremer, M. (2008). Outside funding and the dynamics of participation in community associations. *American Journal of Political Science*, 585–602.
- Horowitz, D. L. (1985). *Ethnic groups in conflict*. University of California Press.
- Kasara, K. (2017). Does local ethnic segregation lead to violence? Evidence from Kenya. *Quarterly Journal of Political Science*, 11(4), 441–470.
- Knox, D., Yamamoto, T., Baum, M., & Berinsky, A. (2019). Design, Identification, and Sensitivity Analysis for Patient Preference Trials. *Journal of the American Statistical Association*, 29, 1–27.
- Miguel, E., Satyanath, S., & Sergenti, E. (2004). Economic shocks and civil conflict: An instrumental variables approach. *Journal of political Economy*, 112(4), 725–753.
- Olken, B. A. (2007). Monitoring corruption: Evidence from a field experiment in Indonesia. *Journal of Political Economy*, 115, 200–249.
- Pearl, J. (2009). *Causality*. Cambridge University Press.
- Pearl, J., & Bareinboim, E. (2014). External validity: From do-calculus to transportability across populations. *Statistical Science*, 29(4), 579–595.
- Scacco, A., & Warren, S. (2018). Can social contact reduce prejudice and discrimination? Evidence from a field experiment in Nigeria. *American Political Science Review*, 112(3), 654–677.
- Wilkinson, S. I. (2004). *Votes and violence: Electoral competition and ethnic riots in India*. Cambridge University Press.